


AGNUS LLM: Robust and Flexible Entity Disambiguation with decoder-only Language Models

Kristian Noullet, Ayoub Ourgani, Niklas Thomas Lakner, Lukas Kinder, Tobias Käfer.

Karlsruhe Institute for Technology (KIT)
Kaiserstraße 12, 76131 Karlsruhe, Karlsruhe, Germany.

Correspondence: noullet@kit.edu

Abstract

Entity disambiguation (ED) links ambiguous mentions in text to entries in a knowledge base and is a core task in entity linking systems. While pretrained decoder-only language models (DLMs) offer strong generalization capabilities, their effective use in ED has been restricted due to sensitivity to candidate order, susceptibility to hallucinated outputs, and potential dataset leakage. We introduce AGNUS  a zero-shot ED framework that addresses these challenges through three core innovations: (1) **order-invariant candidate encoding** via shared positional embeddings and modified autoregressive attention masking, which eliminates bias on input ordering; (2) **constrained decoding** that ensures outputs are restricted to valid candidates, effectively preventing hallucinations; and (3) **synthetic dataset** creation approach as a diagnostic tool for data contamination detection and counteraction. AGNUS eliminates up to 15.2% of F1 variability caused by candidate permutations, delivering consistent and order-robust predictions previously unattainable with autoregressive architectures. In our experiments, AGNUS achieves state-of-the-art performance on four standard ED benchmarks, surpassing prior zero-shot approaches by an average 3.7%. We release code, data including candidate sets, and a synthetic benchmark to support reproducibility and controlled evaluation¹.


1 Introduction

Entity Disambiguation (ED) represents the task of linking ambiguous mentions in text to the correct entity from a provided candidate set and is a core component of knowledge-intensive NLP applications such as question answering, semantic search, and entity linking. While Large Language Models (LLMs) have demonstrated remarkable generalization across diverse tasks, their robust application to

ED remains challenging, particularly in zero-shot settings.

In this work, we identify and address three fundamental limitations of Decoder-Only Language Models (DLMs) when applied to ED:

First, these models are highly sensitive to the order in which candidate entities are presented. Autoregressive generation induces positional bias, leading to substantial prediction variability across permutations of candidate inputs - up to 15.2% (F1) in our experiments. Second, DLMs may produce entities that are not part of the candidate set, undermining system reliability in constrained ED settings. Third, meaningfully evaluating LLMs on existing benchmarks is complicated by potential training data contamination, given the opacity of pretraining corpora.

To overcome these challenges, we propose AGNUS , a robust, zero-shot ED framework for DLMs that can be applied to open-weight models out-of-the-box, requiring no fine-tuning or re-training. AGNUS incorporates two key components to achieve order-robustness and hallucination resistance: (1) **Masked Attention Candidate Set (MACS)** to enforce candidate order-invariant encoding; and (2) **Agnus Contextual Decoding (ACDC)** to restrict decoding to valid candidate entities.

Our approach ensures that candidate entities are represented indistinguishably in terms of position (positional embeddings) and interdependencies (attention) to the underlying model, eliminating the influence of input order on token predictions. At the same time, constrained decoding removes hallucinated outputs without sacrificing the model's contextual reasoning ability.

To audit dataset leakage, we design a synthetic dataset construction methodology for ED. The resulting synthetic dataset serves as a diagnostic tool to detect contamination, evaluate with a potentially lesser degree of contamination and therewith al-

¹<https://github.com/kmdn/agn-dis>

lows us to test model generalization capabilities in a controlled setting.

Our contributions are as follows:

- We propose AGNUS 🧠, a zero-shot ED framework combining:
 - **MACS** for order-invariant candidate encoding.
 - **ACDC** for constrained autoregressive tree-based decoding.
- We introduce a synthetic dataset construction method to gauge benchmark contamination and apply it to the AIDA (Yosef et al., 2011) benchmark.
- We release all code, data including entity candidates and evaluations to support reproducibility² and future comparability.

Across four standard ED benchmarks, AGNUS 🧠 achieves state-of-the-art performance in zero-shot settings, while delivering stable predictions under candidate permutations and eliminating hallucinated outputs.

The remainder of this paper is structured as follows. Section 2 reviews related work in entity disambiguation and recent advances in large language models, with particular attention to dataset contamination, hallucination, and order sensitivity. Section 3 introduces our proposed framework, AGNUS 🧠, detailing the disambiguation setup, our order-invariant encoding method, and constrained decoding strategy. Section 4 presents our experimental setup, results across standard benchmarks, and a comprehensive ablation analysis, followed by a study on contamination detection. Finally, Section 5 concludes the paper and outlines directions for future work.

2 Related Work

2.1 Entity Disambiguation

Entity disambiguation (ED) is a critical task in natural language processing and understanding, where the goal is to map ambiguous entity mentions in text to their correct entries in a knowledge base. Current state-of-the-art ED and entity linking models (van Hulst et al., 2020; Barba et al., 2022; Ayoola et al., 2022; Shavarani and Sarkar,

2023; Xiao et al., 2023b; Ding et al., 2024a; Orlando et al., 2024) make use of various deep learning architectures to outperform more traditional works. In recent years, transformer-based systems, such as BLINK (Wu et al., 2020), REL (van Hulst et al., 2020), SpEL (Shavarani and Sarkar, 2023), DeepType (Raiman and Raiman, 2018) and GENRE (Cao et al., 2021) have taken over the stage with many basing themselves on BERT (Devlin et al., 2019) embeddings. In recent years, LLM-based systems have entered the space with (Sun et al., 2023), (Wang et al., 2023a), (Xiao et al., 2023a), EntGPT (Ding et al., 2024a), ChatEL (Ding et al., 2024b), LLMAEL (Xin et al., 2024) and (Tasawong et al., 2024). Particularly, in (Ding et al., 2024a; Xin et al., 2024; Liu et al., 2024; Vollmers et al., 2025) authors improve LLM-based entity disambiguation by tuning inputs and otherwise providing LLM backbones with context-relevant data.

2.2 Large Language Models

Applying LLMs to ED is accompanied by a multitude of considerations when contrasted with more traditional ED. Among these, there exist benchmark contamination (Section 2.2.1), hallucinations (Section 2.2.2), decoding mechanisms (Section 2.2.3) and order-specific biases (Section 2.2.4) that endanger robust disambiguation. In the following, we address these areas of prior work.

2.2.1 Dataset Contamination

Benchmark contamination in LLMs (Xu et al., 2024) has become a critical issue as models trained on vast amounts of publicly available data may inadvertently ‘memorize’ aspects of popular benchmark datasets, potentially leading to inflated estimates of their true capabilities.

To address these challenges, researchers have started developing various countermeasures (Chen et al., 2025), including dynamic evaluation benchmarks (Wang et al., 2025; Zhu et al., 2024a,b) to effectively prevent pre-benchmarking disclosure. Another measure is to provide a means of evaluation for the degree of contamination (Xu et al., 2024) by computing perplexity (Li, 2023) – by applying the exponential function to the average negative log likelihood over a particular sequence of text to measure a model’s ‘surprise’ (or inverse confidence) for a particular output.

²<https://github.com/kmdn/agn-dis>

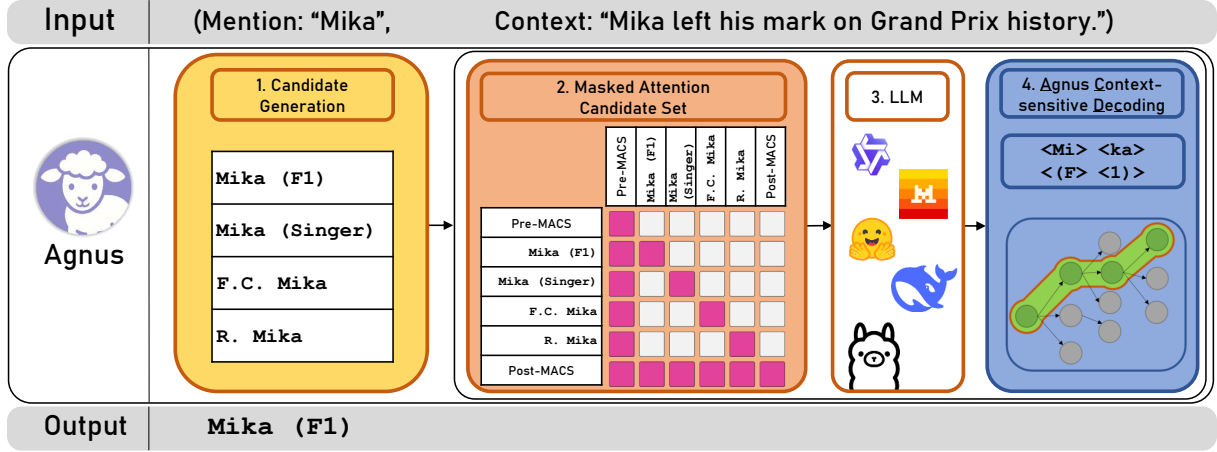


Figure 1: AGNUS Overview – Takes an input document, (1.) generates candidate entities for mentions (e.g. MIKA) using a pre-existing candidate generation method, (2.) applies masked attention and altered positional embeddings to the candidate entity collection (MACS, Section 3.2) and (3.) passes representation to a specified LLM, followed by (4.) constrained decoding (ACDC, Section 3.3) for context-sensitive disambiguation and returns the disambiguated entity (e.g. Mika (F1)).

2.2.2 Hallucinations

Despite remarkable capabilities in generating human-like text, LLMs may produce factual inaccuracies or nonsensical sequences, a phenomenon referred to as *hallucination* (Huang et al., 2025). The underlying causes of hallucinations are an active area of research. Some potential contributing factors include the vast scale of the training data, containing potentially noisy data (Petroni et al., 2021; Ji et al., 2023) and the autoregressive nature of text generation based on prior tokens (Holtzman et al., 2020; Maynez et al., 2020). The presence of hallucinations poses a significant challenge for the reliable application of LLMs on downstream Natural Language Processing (NLP) tasks, posing issue for robust and trustworthy ED. Recent research efforts have started counteracting hallucinations through retrieval augmentation, fact verification and the incorporation of knowledge graphs (Lewis et al., 2020; Pusch and Conrad, 2024).

In this paper, we eliminate the possibility for entity candidate hallucinations by defining a specialised constrained decoding strategy for ED.

2.2.3 Constrained Decoding

Early work on LLMs (Brown et al., 2020; Radford et al., 2019) demonstrated that decoder-only language models process natural language prompts effectively without an enforced schema, meaning that input-output pairs are structurally not bound by predefined templates or grammars. This flexibility allows for broad applicability but introduces challenges in reliability, consistency, and controlla-

bility (Bender et al., 2021).


To mitigate challenges of unstructured interaction, researchers have developed various prompt engineering methods (Sahoo et al., 2024; Ouyang et al., 2022a; Madaan et al., 2023; Wei et al., 2022) to implicitly guide, but not force DLMs towards more structured outputs. Therefore, constrained decoding (Beurer-Kellner et al., 2024) approaches to enforce strict restrictions on LLM text generation have been developed and even started being applied to the domain of entity linking (Vollmers et al., 2025).

2.2.4 Order Bias

Prior work has established that modern generative large language models demonstrate inherent tendencies toward positional preferences when processing ordered lists of candidate answers (Pezeshkpour and Hruschka, 2023; Wei et al., 2024; Zheng et al., 2023; Kinder et al., 2025) and that these are also sensitive towards the arrangement order of otherwise identical answer collections (Dominguez-Olmedo et al., 2023; Li et al., 2023; Li and Gao, 2024; Wang et al., 2023b, 2024a; Xue et al., 2024). Approaches to mitigation include compensation for positional preferences (Wei et al., 2024; Zhao et al., 2021), systematic permutation averaging and applying multiple forward passes with varied option sequences (Pezeshkpour and Hruschka, 2023; Wang et al., 2023b), as well as reasoning-enhanced strategies (Wang et al., 2024a,b) to attenuate sequence dependence. While we alleviate certain biases, we note that AGNUS

may still suffer from other forms, e.g. domain as found in (Noullet et al., 2025). AGNUS employs a method to mitigate candidate order bias without requiring additional training by adapting the approach from (Kinder et al., 2025) – which investigated the effects of altered positional embeddings on generated output – to entity disambiguation.

3 AGNUS

In this section we introduce AGNUS , our proposed approach for LLM-based robust entity disambiguation. In Figure 1, we present AGNUS: from (*Step 1.*) generating entity candidates for each mention using the DBpedia Lookup³ service – chosen for improved reproducibility; (*Step 2.*) applying combined masked attention and position-specific shared positional embeddings (*Masked Attention Candidate Set*, Section 3.2) based on (Kinder et al., 2025); (*Step 3.*) passing the encoded inputs to a chosen model; (*Step 4.*) constrained decoding (ACDC) to final disambiguation for the input document "*Mika left his mark on Grand Prix history.*" and entity mention *Mika*, yielding contextually disambiguated entity *Mika* (F1)⁴.

3.1 Disambiguation Setup

AGNUS represents an approach leveraging DLMs for the task of disambiguating entities based on entity candidate information while mitigating DLM-specific challenges. AGNUS is built on top of the *Combining Linking Techniques* (Noullet et al., 2023) framework and acts as a knowledge-base agnostic entity disambiguator. For disambiguation, AGNUS takes as input a document providing context, a mention and a collection of candidate entities generated via pre-existing candidate generation approaches.

Due to leveraging the contextual disambiguation capabilities of DLMs, AGNUS does not require candidate entities to solely be a knowledge base-backed IRI⁵. Instead, candidate entity representation may additionally take any identifying or meaningful form, such as a description, label, type or combination thereof. For each mention contained within an input document, we generate a fixed candidate set, employing candidates generated with DBpedia Lookup⁶. Each candidate collection is

encoded using MACS (Section 3.2). Subsequently, the resulting encoded prompt is transmitted as a whole to the underlying DLM for contextual parsing and decoded via ACDC (Section 3.3).

3.2 Masked Attention Candidate Set

Text sequences encoded on modern generative language models rely on underlying position-influenced attention mechanisms and positional embeddings to signal the order of token appearance within a sequence (Kinder et al., 2025). This affects desired order-invariant sequences, such as candidate collections – an undesirable property for entity disambiguation. To render an LLM order-agnostic for parts of a sequence, we tackle both aspects: modify positional embeddings (Section 3.2.1) for candidate entities to simulate similar positions and mask the attention mechanism between entity candidates (Section 3.2.2) to the underlying language model.

Each candidate collection is encoded using MACS, embedded into its original textual encoding with text preceding (Pre-MACS) and succeeding (Post-MACS) the collection being encoded in standard LLM-specific fashion.

3.2.1 Positional Embedding

Every sequence of tokens is attributed a certain range of positional embedding values within its LLM-encoded representation. Within a MACS-encoded collection, every token making up an entity candidate is modified to appear as sharing a similar range of positions (see visualization Fig. 2) as other candidates to the underlying LLM.

To do so, we define relative position $i \in [0, \dots, n_{c_j} - 1]$ of each token $t_{c_j,i}$ for entity candidate representation $c_j \in C$ s.t. n_{c_j} is the number of tokens for entity c_j and collection of all candidate entities C for a given mention and T_{c_j} the set of all tokens for c_j : $\forall t_{c_j,i} : i \in [0, \dots, \max_{c \in C} (|n_c - 1|)]$.

Therefore as visualised in Figure 2, the shared range of possible positional embeddings is defined by the token-wise longest candidate within a MACS collection and starts for each candidate at the end of prior sequence's token (PRE-MACS) and afterwards continues the candidate encoding with the succeeding sequence's (POST-MACS) first positional embedding.

³<https://lookup.dbpedia.org/>

⁴https://en.wikipedia.org/wiki/Mika_Häkkinen

⁵<https://wikipedia.org/wiki/>

Internationalized_Resource_Identifier

⁶<https://github.com/dbpedia/dbpedia-lookup>

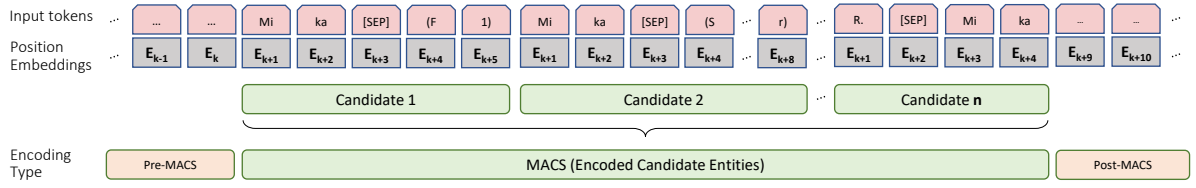


Figure 2: MACS – Positional Embeddings: Each candidate entity entry is encoded as being on the same range of positions for the length of their contents. Candidate entity entries’ positional identifiers are **shared** across common lengths and encoded analogously. Post-MACS – any tokens after a MACS block – starting positional embedding is computed as being subsequent to the longest option contained within MACS (candidate) entries.

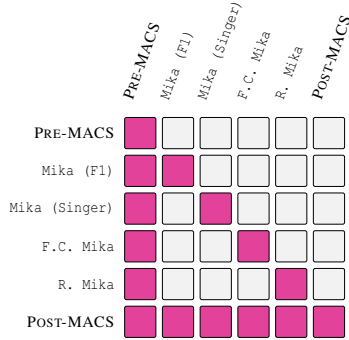


Figure 3: MACS – Causal mask: Example from Fig. 1 for entity candidate representations for entities "Mika (F1)", "Mika (Singer)", "F.C. Mika", "R. Mika". Grey cells signify blocked attention whereas pink signifies enabled attention. Intra entity attention and attention from tokens preceding (Pre-MACS) and succeeding (Post-MACS) MACS is preserved normally s.t. subsequent tokens attend to prior ones.

3.2.2 Causal Mask

To encode a collection of entity candidates in an order-invariant fashion to the underlying LLM, we apply an adapted version of the commonly-employed triangular attention matrix as causal mask (see Figure 3). Entities within a collection cannot attend to one another (grey entries), but do attend (pink entries) – and are attended to – in otherwise usual LLM fashion to their own prior tokens (diagonal entries) and rest of the token sequence (to PRE-MACS and by POST-MACS). This means that tokens within each candidate’s representation continue attending to each other, such as `<ka>` to `<Mi>` within our example.


3.3 Agnus Contextual Decoding

Generative LLMs may hallucinate information in unexpected fashions. This ranges from a corrupt expected result format to non-existing options. Due to the nature of entity disambiguation, only given options may be produced. As such, we define an input-flexible grammar based on entity candidates.

We implement this grammar in a logits processor⁷ that filters forbidden tokens at each generation step until a single disambiguated entity remains.

Formally: Let the set of candidate sequences be $O = \{o_1, \dots, o_n\}$ where each candidate option $o_i \in \Sigma^{l_i}$ is a sequence of length l_i . The vocabulary is defined as $\Sigma = \{t_k^i \mid i \in \{1, \dots, n\}, k \in \{1, \dots, l_i\}\} \cup \{\text{EOS}\}$. We then define the set of nonterminals as $V = \{X_i^k \mid i \in \{1, \dots, n\}, k \in \{0, \dots, l_i\}\}$ where X_i^k denotes the state after generating the first k tokens of candidate c_i . The start symbol transitions to the initial state of each candidate: $S \rightarrow X_1^0 \mid X_2^0 \mid \dots \mid X_n^0$. For each o_i , we define the following transitions: $X_i^k \rightarrow t_{k+1}^i X_i^{k+1}, \forall k \in \{0, \dots, l_i - 1\}$, $X_i^{l_i} \rightarrow \text{EOS}$.

4 Experiments and Results

AGNUS  combines techniques to create an LLM-enabled approach to robust entity candidate disambiguation. In this section, we conduct experiments to evaluate AGNUS with different configurations regarding representations of entity candidates, LLMs, our candidate encoding (MACS) and our constrained decoding (ACDC). We report entity disambiguation results in comparison to prior work in Table 1.

4.1 Technical Details

All our experiments were run on a server with NVIDIA RTX 4090 (24GB vRAM), 1TB RAM, 128 CPU cores, Debian (Bookworm), CUDA 12.5 and Python 3.11. As for LLMs, we decided on instruct models for our experiments such that they would run on our hardware and be comparable in size, leading to the following selection: Mistral (7B-Instruct) (Jiang et al., 2023), Llama2 (7B) (Touvron et al., 2023), Llama3 (8B-Instruct) (Dubey et al., 2024) and Qwen (2.5-7B-

⁷<https://anonymous.4open.science/r/Agnus/src/agnus/pipeline/llm.py#L356>

Instruct) (Yang et al., 2024) – for the rest of the paper we omit detailed version specifications.

4.2 Evaluation

In our experiments, we outperform related work on 4 out of 5 common datasets (AIDA (Yosef et al., 2011), KORE 50^{DYWC} (Noullet et al., 2020), MSNBC (Cucerzan, 2007), ACE04 (Ratinov et al., 2011), AQUAINT (Milne and Witten, 2008)) in zero-shot settings despite our underlying LLMs’ relatively modest parameter count⁸.

We report our ED F1 results in Table 1 – the *top block* lists scores for trained or finetuned approaches, the *bottom block* compares zero-shot methods. Our model performs strongly across most datasets and even surpasses finetuned or trained prior work in certain cases. Despite being a zero-shot approach, AGNUS attains overall new state-of-the-art results for KORE 50 (82.3%) and ACE04 (95.5%). Unsurprisingly when evaluating on AIDA, approaches trained on AIDA outperform ours, but AGNUS (86.7%) exceeds second-ranked zero-shot approach EntGPT-P (Ding et al., 2024a) (82.1%) F1 measure by 4.6%. Evaluating against KORE 50, AGNUS reaches 82.3% in comparison to ChatEL’s 78.7%, surpassing it by 3.6%. As for ACE and AQUAINT, our results (95.5% and 87.5%) improve upon EntGPT-P’s (91.8% and 79.1%) respectively by 3.7% and 8.4%. For MSNBC, we do not beat the state-of-the-art for zero-shot entity disambiguation and instead reach 82.4%, underperforming ChatEL (Ding et al., 2024b) (88.1%) by 5.7% and finetuned state-of-the-art CoherentED (96.3%) by 13.9%. While AGNUS yields improvements across some benchmarks, we consider our primary benefit lying in enhancing disambiguation robustness via order invariance for candidates and by preventing structurally invalid outputs.

4.3 Ablation Study

AGNUS employs multiple techniques to mitigate issues relating to LLM-based ED. Particularly, AG-

NUS relies on LLMs for disambiguation, MACS for order-invariant candidate encoding, ACDC for entity decoding and particularly candidates’ representation. In our ablation study, we therefore design experiments to verify the impact of these aspects on model results by investigating candidate representation (Section 4.3.1), LLM selection (Section 4.3.2), MACS (Section 4.3.3) and ACDC (Section 4.3.4).

4.3.1 Candidate Representation

To validate LLM disambiguation capabilities based on contextual candidate entity information, we apply AGNUS to candidate representations of different entity information types. We selected DBpedia (& Wikipedia) entity IRIs, entity types, textual entity descriptions and labels as meaningful entity information characterising entity candidates for our experiments (Tables 2 and 3). We note that in Table 2 across all datasets, IRI-based representations perform best with an average F1 performance of 86.9%, outperforming labels by 10.2% – with a tie of 87.9% for AQUAINT. For all datasets beside KORE 50 and AQUAINT, descriptions reach the second-highest score (avg.: 73.5%), but are still surpassed by labels (76.7%) on average by 3.2%. We note that the shorter and more unique a representation is, the better AGNUS seems to perform. In our experiments, we find effects of representation depend on benchmarked dataset with representation-based score differences ranging from 5.6% (ACE04) to 34.9% (AQUAINT) with a mean of 21.12% across our 6 datasets.

4.3.2 Large Language Model

To verify our approach’s generalizability across LLMs, we run AGNUS on 4 LLMs: Llama2 (Touvron et al., 2023), Llama3 (Dubey et al., 2024), Qwen (Yang et al., 2024) and Mistral (Jiang et al., 2023). In Table 3, we notice similar trends across most LLMs for the AIDA dataset with Llama2 representing a slight outlier: All other LLMs attain respective top results using IRIs (Qwen: 84.6%, Mistral: 86.7%, Llama3: 84.0%) as candidate information, whereas our outlier LLM manages to slightly improve on its 80.9% F1 measure, reaching 81.3% by employing labels as candidate representation. Typically, Llama3, Mistral and Qwen reach similar results to each other using IRIs (84.0% – 86.7%) and descriptions (70.2% – 76.3%) as candidate representations. Using labels, Qwen plummets down to 64.6%, whereas Llama3 (74.5%) and Mis-

⁸We note that our models are at least an order of magnitude smaller and that Ding et al. (2024b) argue model parameter count having a significant influence on the entity disambiguation task. EntGPT (Ding et al., 2024a) and ChatEL (Ding et al., 2024b) employ Llama2 70B (Touvron et al., 2023) and GPT-3.5 (Ouyang et al., 2022b); ChatEL (Ding et al., 2024b) additionally makes use of PaLM 540B (Chowdhery et al., 2023) and GPT-4 (OpenAI, 2023). OpenAI has not disclosed parameter counts for GPT-3.5 and GPT-4, but each of them is assumed to have at least 175B parameters, with rumors claiming GPT-4 having 1.76 trillion parameters according to <https://en.wikipedia.org/wiki/GPT-4>.

Trained (or finetuned) for ED						
Model	AIDA	KORE 50	MSNBC	ACE04	AQUAINT	Mean
End2End (Kolitsas et al., 2018)	0.891	0.569	0.933	0.892	0.894	0.836
GENRE (Cao et al., 2021)	0.933	0.542	0.943	0.901	0.899	0.844
REL (van Hulst et al., 2020)	0.928	<u>0.618</u>	0.935	0.897	0.873	0.850
ReFinED (Ayoola et al., 2022)	<u>0.939</u>	0.567	0.941	0.908	0.918	0.855
LLMAEL × ReFinED _{FT} (Xin et al., 2024)	0.923	-	0.888	0.881	0.891	0.900
EntGPT-I (GPT3.5) (Ding et al., 2024a)	0.920	0.753	0.922	0.937	0.906	0.888
EOEDbMSL (Tasawong et al., 2024)	0.941	-	0.935	0.917	0.894	0.922
ExtEnD (Barba et al., 2022)	0.926	-	<u>0.947</u>	0.918	<u>0.916</u>	<u>0.927</u>
CoherentED (Xiao et al., 2023b)	0.894	-	0.963	<u>0.934</u>	0.946	0.934

LLM 0-shot ED						
Model	AIDA	KORE 50	MSNBC	ACE04	AQUAINT	Mean
ChatEL (Ding et al., 2024b)	-	<u>0.787</u>	0.881	0.893	0.767	<u>0.832</u>
EntGPT-P (GPT3.5) (Ding et al., 2024a)	<u>0.821</u>	0.716	<u>0.867</u>	<u>0.918</u>	<u>0.791</u>	0.823
EntGPT-P (Llama2 70B) (Ding et al., 2024a)	0.708	0.647	0.741	0.746	0.635	0.695
Ours – AGNUS (Llama2 8B)	0.809	0.529	0.562	0.897	0.576	0.675
Ours – AGNUS (Mistral)	0.867	0.823	0.824	0.955	0.875	0.869
Baseline: Mistral (hidden candidates)	0.791	0.794	0.739	0.953	0.720	0.799
Ablation: w.o. MACS (best)	0.865	0.811	0.814	0.962	0.907	(0.872)
Ablation: w.o. MACS (worst)	0.833	0.779	0.766	0.950	0.847	(0.835)

Table 1: ED evaluation table – **Upper category**: ED systems trained or finetuned for ED (mainly with AIDA). **Lower category**: 0-shot ED systems. Top scores per column and category **bolded**, second highest underlined. Scores obtained from respective papers. Note that baseline with *hidden candidates* also uses matching to candidates (else naive results would tend to 0) and MACS ablations are run over multiple iterations, showing score variability.

Entity Representation	AIDA	AIDA-Syn	KORE 50	MSNBC	ACE04	AQUAINT	Mean
AGNUS (w. IRI)	0.867	0.863	0.823	0.824	0.955	0.879	0.869
AGNUS (w. Label)	0.743	0.706	<u>0.785</u>	0.589	0.899	0.879	0.767
AGNUS (w. Type)	0.705	0.719	0.595	0.591	0.934	0.530	0.679
AGNUS (w. Description)	<u>0.763</u>	<u>0.790</u>	0.515	<u>0.679</u>	<u>0.954</u>	<u>0.706</u>	0.735
Mean	0.769	0.770	0.679	0.671	0.936	0.748	0.762

Table 2: Ablation Study (Candidate Representation over datasets): AGNUS (Mistral) F1 measures on AIDA, AIDA-Syn, KORE 50, MSNBC, ACE04 and AQUAINT with different candidate entity representations (IRI, label, entity type, entity description), along with per representation and per dataset averages. Top entry by dataset in **bold**, second underlined.

Model	AIDA				Mean
	IRI	Label	Type	Description	
AGNUS (Qwen)	0.846	0.646	0.422	0.721	0.659
AGNUS (Mistral)	0.867	0.743	0.705	0.763	0.770
AGNUS (Llama2)	0.809	0.813	0.560	0.565	0.687
AGNUS (Llama3)	0.840	0.745	0.392	0.702	0.670
AGNUS (Llama3) w.o. ACDC	0.765	0.698	0.331	0.677	0.618

Table 3: Ablation Study (LLM, Candidate Representation, ACDC): AGNUS F1 measures for different types of candidate representations for Qwen, Mistral, Llama2, Llama3 and without constrained decoding via ACDC. AGNUS without ACDC utilises fuzzy search, ranking reply and candidate, matching to candidate with highest similarity.

tral (74.3%) attain F1 scores close to each other. For type candidate information, Mistral (70.5%) noticeably outperforms Qwen (42.2%) and Llama3 (39.2%); Llama2 manages to outperform its successor Llama3 (56.0%). Llama3 (70.2%) and Qwen

(72.1%) handle descriptions as meaningful entity information comparably well with Mistral (76.3%) performing slightly better and Llama2 (56.5%) displaying worst results.

4.3.3 Masked Attention Candidate Set

We investigate how MACS affects qualitative results and whether it actually renders disambiguation order-invariant. To this end, we run experiments shuffling candidates over 10 iterations and display results in Figure 4. Our experiments over 3 different LLMs (Llama3, Mistral, Qwen) display how disambiguation varies without the use of MACS and remains unchanged when applying MACS. Order invariance persists across all 10 iterations of shuffled candidates when MACS is employed whereas not applying the causal mask to candidate entities yields result variations. With-

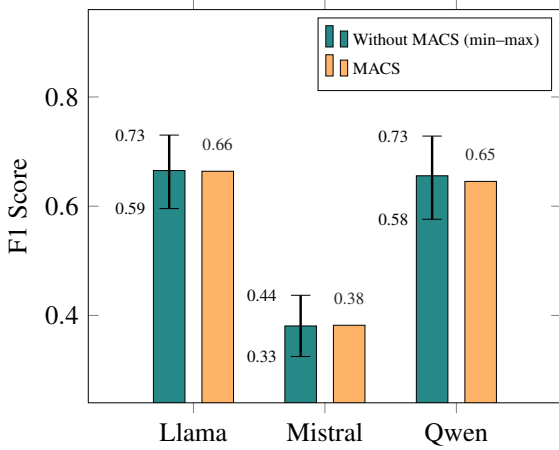


Figure 4: Ablation Study (MACS) - F1 Score Variability: Error Bar plot on disambiguation without (*left*) and with MACS (*right*) with randomised candidate shuffles over 10 iterations with Llama3, Mistral and Qwen on perplexity decoding – disambiguates to highest confidence – for AIDA. MACS and non-MACS results are similar on average. Without MACS, performance varies (Llama3: 13.5%, Mistral: 11.2%, Qwen: 15.2%).

out MACS, Llama3 averages at 66.53% (MACS: 66.40%) and varies between 59.56% – 73.02%, a difference of 13.46%. Mistral on the other hand varies in the range of 32.47% – 43.69%, averaging at 38.07% without MACS across iterations of candidate shuffles (with MACS: 38.20%). In Table 1, we also display MACS ablations over 2 iterations, from one non-MACS execution to another exhibiting 3.7% F1 difference on average and beating AGNUS in ACE04 (96.2% vs. 95.5%). Finally, Qwen also exhibits changes resulting from candidate order changes: with an average of 65.57% (MACS: 64.54%) its candidate order-dependant results vary within the range 57.61% – 72.84%. Based on our experiments, we conclude that MACS effectively removes order-based bias from candidates with an overall minor average reduction in F1 score.

4.3.4 Agnus Contextual Decoding

In Table 3, additionally to checking out the impact of candidate representations across language models, we also evaluate AGNUS without our constrained decoding method (ACDC). We find that AGNUS hallucinates across the board, decreasing F1 for all types of candidate representation. In non-ACDC experiments, we apply fuzzy matching to improve the likelihood of finding at least one entity. Exact disambiguation to candidate matches in our zero-shot experiments yield extremely subpar results (close to 0) and would otherwise be misrepresenting

the added value of our robustness-oriented approach. On average, F1 performance without ACDC is lowered by 5.2%, the largest drops appearing with IRI (-7.5%) and type (-6.1%) candidate representations, followed by label (-4.7%) and descriptions (-2.5%).

4.4 Contamination Detection

To diagnose potential contamination, we employ perplexity (Li, 2023) to quantify a model’s uncertainty for a given token sequence prediction. Perplexity reflects the inverse likelihood assigned to a particular token sequence by a model: lower perplexity indicates higher predictive confidence and a higher likelihood of contamination. To detect contamination and evaluate the generalizability of DLMs, we propose synthetically generating a novel dataset derived from an existing one by replacing each entity mention with a distinct, contextually similar mention and corresponding entity. We apply our method with the DeepSeek-R1 (DeepSeek-AI et al., 2024) model⁹ to AIDA (Yosef et al., 2011) and release AIDA-Syn¹⁰. For each sequence, we produced five mention-entity options, but for AIDA-Syn only one was retained to reduce the risk of future pretraining exposure. All alternatives, along with a generation script, are made available¹¹. To assess contamination levels across different LLMs, we introduce a modified decoding strategy, illustrated in Figure 4 with disambiguation performed by selecting a candidate entity with highest confidence. The model that performs worst with this strategy is presumed to be least contaminated. Our findings show that Mistral (Jiang et al., 2023) yields the lowest performance with a perplexity-based decoding method on AIDA, suggesting being least affected by benchmark contamination. Applying the same decoding strategy with AIDA-Syn, F1 score decreases from 38.20% to 22.82%, a substantial relative drop of 15.38%. This reduction supports the hypothesis that AIDA-Syn exhibits reduced contamination and that our underlying DLMs may suffer from contamination.

⁹Version from May 2025: <https://www.deepseek.com/>

¹⁰Made up of 888 documents. More details in Appendix.

¹¹<https://github.com/kmdn/agn-dis>

5 Conclusion

We propose a set of techniques to enable robust LLM-based entity disambiguation by addressing the issues of unwanted order bias and hallucinations in an entity disambiguation component we dub AGNUS 🐺. Our experimental results show that our zero-shot approach outperforms prior work on average by 3.5%. Further, we introduce a methodology to detect data contamination and publish a novel diagnostic dataset AIDA-Syn. Advantages of our approach to the area of ED are that contextual information may be utilised through the application of LLMs of choice, enabling an ease of candidate representation definitions. This could be particularly beneficial in the case of incomplete, imperfect or not well-connected knowledge bases. While our approach using MACS and ACDC yields modest improvements across benchmarks, our primary benefit lies in enhancing output robustness and controlling generation behavior, particularly in cases where unconstrained and order-variant decoding leads to semantically or structurally invalid outputs. The only requirement we have for entity disambiguation is to have meaningful textual candidate representations. Therefore, AGNUS 🐺 can be used for multi-Knowledge Graph (KG), KG-agnostic fashions and novel entity settings – provided entity-defining information is present.

In future work, we plan to improve upon our research by finetuning models, testing scalability with large numbers of candidates, designing novel decoding strategies and testing our approach over other knowledge bases.

6 Limitations

Due to our introduction of order-invariance by application of a causal mask, modifying positional IDs and introducing a custom logits processor, we are limited to open-weight DLMs, making evaluation with DeepSeek (DeepSeek-AI et al., 2024), GPT-3.5 (Ouyang et al., 2022b), GPT-4 (OpenAI, 2023) impossible, unfortunately. Also, due to hardware restrictions, we are limited to evaluating on significantly smaller models than related work – with us running experiments on models with around 10 billion parameters. For instance, in (Ding et al., 2024a), authors employ GPT3.5 and Llama, claiming that results improve with increased model size.

AGNUS counteracts candidate order bias, but still suffers from other biases, such as domain (Noullet et al., 2025) or data bias linked to

the underlying LLM.

Fundamentally, we design a causal mask due to being interested in disambiguating entities with causal decoder-only language models. The general idea could likely be transferred to other types of language models by analogously adapting the causal mask to fit another paradigm’s attention masking strategy.

Alike other deep learning approaches to entity disambiguation, AGNUS is limited by its generated candidate sets and by only working with candidate entities that have some form of textual label, description, types or otherwise meaningful information for a LLM to predict.

Our prompt design (for further details, see Figure 6) does not take into account character offsets within input documents due to DLMs running into issues when handling numerical values. Consequently, AGNUS may run into issues when the same mention at different positions in a document refers to different entities. (Example: Tim, CEO of **Apple**, likes to eat an **apple** a day.)

While ACDC does mitigate hallucinations, a given LLM’s next token prediction may be to continue with non-entity tokens, such as a greeting, acknowledgment of task or similar, therewith potentially negatively affecting entity disambiguation depending on decoding algorithm. Designing a specific decoding strategy to include such behaviour could prove to be a benefit in future endeavours.

In this paper, our underlying models are not finetuned for the entity disambiguation task nor given particular domain-specific information that could boost context and potential results. Therefore, we concede that going for a few-shot approach could yield improved results.

Further, despite having the out-of-the-box structural capabilities for it, we could not evaluate our approach on knowledge bases other than Wikipedia and DBpedia due to not being aware of comparable and valid evaluation benchmarks for it.

DLMs are language-dependant and have mainly been trained with English in mind. Therefore, results may vary greatly when our approach is applied to other languages.

Regarding evaluation contamination and the creation of AIDA-Syn, we did not go as in-depth explaining our procedure, safeguards against LLM hallucinations, inherent surrounding bias as we would have liked, nor provide in-depth statistics or analyses. We introduce it mainly as a diagnostic tool to evaluate our approach and show that

despite there being novel entities and candidates, AGNUS is capable of attaining similar results as for the non-synthetic version with the suggested least contaminated LLM.

References

- Tom Ayoola, Shubhi Tyagi, Joseph Fisher, Christos Christodoulopoulos, and Andrea Pierleoni. 2022. [Re-fined: An efficient zero-shot-capable approach to end-to-end entity linking](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Track, NAACL 2022, Hybrid: Seattle, Washington, USA + Online, July 10-15, 2022*, pages 209–220. Association for Computational Linguistics.
- Edoardo Barba, Luigi Procopio, and Roberto Navigli. 2022. ExtEnD: Extractive entity disambiguation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, Online and Dublin, Ireland. Association for Computational Linguistics.
- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. [On the dangers of stochastic parrots: Can language models be too big?](#) In *FAccT '21: 2021 ACM Conference on Fairness, Accountability, and Transparency, Virtual Event / Toronto, Canada, March 3-10, 2021*, pages 610–623. ACM.
- Luca Beurer-Kellner, Marc Fischer, and Martin T. Vechev. 2024. [Guiding llms the right way: Fast, non-invasive constrained generation](#). In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. 2021. [Autoregressive entity retrieval](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Simin Chen, Yiming Chen, Zexin Li, Yifan Jiang, Zhongwei Wan, Yixin He, Dezhi Ran, Tianle Gu, Haizhou Li, Tao Xie, and Baishakhi Ray. 2025. [Recent advances in large language model benchmarks against data contamination: From static to dynamic evaluation](#). *CoRR*, abs/2502.17521.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2023. [Palm: Scaling language modeling with pathways](#). *J. Mach. Learn. Res.*, 24:240:1–240:113.
- Silviu Cucerzan. 2007. [Large-scale named entity disambiguation based on wikipedia data](#). In *EMNLP-CoNLL 2007, Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, June 28-30, 2007, Prague, Czech Republic*, pages 708–716. ACL.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, and Wangding Zeng. 2024. [Deepseek-v3 technical report](#). *CoRR*, abs/2412.19437.


- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.
- Yifan Ding, Amrit Poudel, Qingkai Zeng, Tim Weninger, Balaji Veeramani, and Sanmitra Bhat-tacharya. 2024a. [Entgpt: Linking generative large language models with knowledge bases](#). *CoRR*, abs/2402.06738.
- Yifan Ding, Qingkai Zeng, and Tim Weninger. 2024b. [Chatel: Entity linking with chatbots](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy*, pages 3086–3097. ELRA and ICCL.
- Ricardo Dominguez-Olmedo, Moritz Hardt, and Celestine Mendler-Dünner. 2023. Questioning the survey responses of large language models. *arXiv preprint arXiv:2306.07951*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. [The llama 3 herd of models](#). *CoRR*, abs/2407.21783.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The curious case of neural text degeneration](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *ACM Transactions on Information Systems*, 43(2):1–55.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation](#). *ACM Comput. Surv.*, 55(12):248:1–248:38.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *CoRR*, abs/2310.06825.
- Lukas Kinder, Lukas Edman, Alexander Fraser, and Tobias Käfer. 2025. [Positional overload: Positional debiasing and context window extension for large language models using set encoding](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 3896–3908. Association for Computational Linguistics.
- Nikolaos Kolitsas, Octavian-Eugen Ganea, and Thomas Hofmann. 2018. [End-to-end neural entity linking](#). In *Proceedings of the 22nd Conference on Computational Natural Language Learning, CoNLL 2018, Brussels, Belgium, October 31 - November 1, 2018*, pages 519–529. Association for Computational Linguistics.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive NLP tasks](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Ruizhe Li and Yanjun Gao. 2024. Anchored answers: Unravelling positional bias in gpt-2’s multiple-choice questions. *arXiv preprint arXiv:2405.03205*.
- Yucheng Li. 2023. [Estimating contamination via perplexity: Quantifying memorisation in language model evaluation](#). *CoRR*, abs/2309.10677.
- Zongjie Li, Chaozheng Wang, Pingchuan Ma, Daoyuan Wu, Shuai Wang, Cuiyun Gao, and Yang Liu. 2023. Split and merge: Aligning position biases in large

- language model based evaluators. *arXiv preprint arXiv:2310.01432*.
- Xukai Liu, Ye Liu, Kai Zhang, Kehang Wang, Qi Liu, and Enhong Chen. 2024. [Onenet: A fine-tuning free framework for few-shot entity linking via large language model prompting](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 13634–13651. Association for Computational Linguistics.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. [Self-refine: Iterative refinement with self-feedback](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan T. McDonald. 2020. [On faithfulness and factuality in abstractive summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 1906–1919. Association for Computational Linguistics.
- David N. Milne and Ian H. Witten. 2008. [Learning to link with wikipedia](#). In *Proceedings of the 17th ACM Conference on Information and Knowledge Management, CIKM 2008, Napa Valley, California, USA, October 26-30, 2008*, pages 509–518. ACM.
- Andrea Moro, Alessandro Raganato, and Roberto Navigli. 2014. Entity Linking meets Word Sense Disambiguation: a Unified Approach. *Transactions of the Association for Computational Linguistics (TACL)*, 2:231–244.
- Kristian Noullet, Rico Mix, and Michael Färber. 2020. [KORE 50^{DyWC}: An Evaluation data set for Entity Linking based on DBpedia, YAGO, Wikidata, and crunchbase](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 2389–2395, Marseille, France. European Language Resources Association.
- Kristian Noullet, Ayoub Ourgani, and Michael Färber. 2023. [A full-fledged Framework for Combining Entity Linking Systems and Components](#). In *Proceedings of the 12th Knowledge Capture Conference 2023, K-CAP '23*, page 148–156, New York, NY, USA. Association for Computing Machinery.
- Kristian Noullet, Ayoub Ourgani, Niklas Lakner, and Tobias Käfer. 2025. [Linking with Bias: Domain-Specific Behaviour in Entity Linking Systems](#). In *Volume 62: Linking Meaning: Semantic Technologies Shaping the Future of AI*, pages 122–139. IOS Press, Vienna, Austria.
- OpenAI. 2023. [GPT-4 technical report](#). *CoRR*, abs/2303.08774.
- Riccardo Orlando, Pere-Lluís Huguet Cabot, Edoardo Barba, and Roberto Navigli. 2024. Retrieve, read and link: Fast and accurate entity linking and relation extraction on an academic budget. In *Findings of the Association for Computational Linguistics: ACL 2024*, Bangkok, Thailand. Association for Computational Linguistics.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022a. [Training language models to follow instructions with human feedback](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022b. [Training language models to follow instructions with human feedback](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. 2021. [KILT: a benchmark for knowledge intensive language tasks](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 2523–2544. Association for Computational Linguistics.
- Pouya Pezeshkpour and Estevam Hruschka. 2023. Large language models sensitivity to the order of options in multiple-choice questions. *arXiv preprint arXiv:2308.11483*.
- Francesco Piccinno and Paolo Ferragina. 2014. [From tagme to WAT: a new entity annotator](#). In *ERD'14, Proceedings of the First ACM International Workshop on Entity Recognition & Disambiguation, July 11, 2014, Gold Coast, Queensland, Australia*, pages 55–62. ACM.
- Larissa Pusch and Tim O. F. Conrad. 2024. [Combining llms and knowledge graphs to reduce hallucinations in question answering](#). *CoRR*, abs/2409.04181.

- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Jonathan Raiman and Olivier Raiman. 2018. [Deeptype: Multilingual entity linking by neural type system evolution](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 5406–5413. AAAI Press.
- Lev-Arie Ratinov, Dan Roth, Doug Downey, and Mike Anderson. 2011. [Local and global algorithms for disambiguation to wikipedia](#). In *The 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference, 19-24 June, 2011, Portland, Oregon, USA*, pages 1375–1384. The Association for Computer Linguistics.
- Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. [A systematic survey of prompt engineering in large language models: Techniques and applications](#). *CoRR*, abs/2402.07927.
- Hassan Shavarani and Anoop Sarkar. 2023. [Spel: Structured prediction for entity linking](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 11123–11137. Association for Computational Linguistics.
- Chenchen Sun, Yuyuan Jin, Derong Shen, Tiezheng Nie, Xite Wang, and Yingyuan Xiao. 2023. [Enhancing knowledge graph attention by temporal modeling for entity alignment with sparse seeds](#). In *Database Systems for Advanced Applications - 28th International Conference, DASFAA 2023, Tianjin, China, April 17-20, 2023, Proceedings, Part II*, volume 13944 of *Lecture Notes in Computer Science*, pages 639–655. Springer.
- Panuthep Tasawong, Peerat Limkonchotiawat, Potsawee Manakul, Can Udomcharoenchaikit, Ekapol Chuangsuwanich, and Sarana Nutanong. 2024. [Efficient overshadowed entity disambiguation by mitigating shortcut learning](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 15313–15321. Association for Computational Linguistics.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *CoRR*, abs/2307.09288.
- Johannes M. van Hulst, Faegheh Hasibi, Koen Dercksen, Krisztian Balog, and Arjen P. de Vries. 2020. [Rel: An entity linker standing on the shoulders of giants](#). In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '20*, page 2197–2200, New York, NY, USA. Association for Computing Machinery.
- Ruben Verborgh, Michael Röder, Ricardo Usbeck, and Axel-Cyrille Ngonga Ngomo. 2018. [Gerbil – benchmarking named entity recognition and linking consistently](#). *Semant. Web*, 9(5):605–625.
- Daniel Vollmers, Hamada M. Zahera, Diego Mousallem, and Axel-Cyrille Ngonga Ngomo. 2025. [Contextual augmentation for entity linking using large language models](#). In *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 8535–8545. Association for Computational Linguistics.
- Kehang Wang, Qi Liu, Kai Zhang, Ye Liu, Hanqing Tao, Zhenya Huang, and Enhong Chen. 2023a. [Class-dynamic and hierarchy-constrained network for entity linking](#). In *Database Systems for Advanced Applications - 28th International Conference, DASFAA 2023, Tianjin, China, April 17-20, 2023, Proceedings, Part II*, volume 13944 of *Lecture Notes in Computer Science*, pages 622–638. Springer.
- Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023b. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*.
- Siyuan Wang, Zhuohan Long, Zhihao Fan, Xuanjing Huang, and Zhongyu Wei. 2025. [Benchmark self-evolving: A multi-agent framework for dynamic LLM evaluation](#). In *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 3310–3328. Association for Computational Linguistics.
- Xinpeng Wang, Chengzhi Hu, Bolei Ma, Paul Röttger, and Barbara Plank. 2024a. Look at the text: Instruction-tuned language models are more robust

- multiple choice selectors than you think. *arXiv preprint arXiv:2404.08382*.
- Xinpeng Wang, Bolei Ma, Chengzhi Hu, Leon Weber-Genzel, Paul Röttger, Frauke Kreuter, Dirk Hovy, and Barbara Plank. 2024b. "my answer is c": First-token probabilities do not match text answers in instruction-tuned language models. *arXiv preprint arXiv:2402.14499*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Sheng-Lun Wei, Cheng-Kuang Wu, Hen-Hsen Huang, and Hsin-Hsi Chen. 2024. Unveiling selection biases: Exploring order and token sensitivity in large language models. *arXiv preprint arXiv:2406.03009*.
- Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. 2020. [Scalable zero-shot entity linking with dense entity retrieval](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 6397–6407. Association for Computational Linguistics.
- Zilin Xiao, Ming Gong, Jie Wu, Xingyao Zhang, Linjun Shou, and Daxin Jiang. 2023a. [Instructed language models with retrievers are powerful entity linkers](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 2267–2282. Association for Computational Linguistics.
- Zilin Xiao, Linjun Shou, Xingyao Zhang, Jie Wu, Ming Gong, and Daxin Jiang. 2023b. [Coherent entity disambiguation via modeling topic and categorical dependency](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 7480–7492. Association for Computational Linguistics.
- Amy Xin, Yunjia Qi, Zijun Yao, Fangwei Zhu, Kaisheng Zeng, Bin Xu, Lei Hou, and Juanzi Li. 2024. [LL-MAEL: large language models are good context augmenters for entity linking](#). *CoRR*, abs/2407.04020.
- Cheng Xu, Shuhao Guan, Derek Greene, and M. Tahar Kechadi. 2024. [Benchmark data contamination of large language models: A survey](#). *CoRR*, abs/2406.04244.
- Mengge Xue, Zhenyu Hu, Liqun Liu, Kuo Liao, Shuang Li, Honglin Han, Meng Zhao, and Chengguo Yin. 2024. Strengthened symbol binding makes large language models reliable multiple-choice selectors. *arXiv preprint arXiv:2406.01026*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2024. [Qwen2.5 technical report](#). *CoRR*, abs/2412.15115.
- Mohamed Amir Yosef, Johannes Hoffart, Ilaria Bordino, Marc Spaniol, and Gerhard Weikum. 2011. [AIDA: an online tool for accurate disambiguation of named entities in text and tables](#). *Proc. VLDB Endow.*, 4(12):1450–1453.
- Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *International conference on machine learning*, pages 12697–12706. PMLR.
- Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. 2023. Large language models are not robust multiple choice selectors. In *The Twelfth International Conference on Learning Representations*.
- Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. 2024a. [Dyval: Dynamic evaluation of large language models for reasoning tasks](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Kaijie Zhu, Jindong Wang, Qinlin Zhao, Ruochen Xu, and Xing Xie. 2024b. [Dyval 2: Dynamic evaluation of large language models by meta probing agents](#). *CoRR*, abs/2402.14865.

A Appendix

Over the course of researching and developing AG-NUS , we implemented some further aspects that we could not allude to in depth. Here are some supplemental materials about them that might be of interest to fellow researchers.

A.1 Prompt Setup & Candidate Representation

In Figure 6, we illustrate how our prompts are set up: we provide a system prompt, a user prompt including a task definition, specific mention to disambiguate, input document text, desired output representation type and a list of (by default) 10 permutable candidate entities. We note that since we do not provide offsets (due to DLMs handling them questionably), we potentially run into issues when multiple mentions - referring to different entities - are to be disambiguated.

Prompts for different representation types are constructed analogously to the URI-based prompt (Figure 5) from our prompt template (Figure 6) for simple representation types. For combined ones, the process is similar with prompt template seen in Figure 7, but utilises 2 (different) representations instead.

A.1.1 IRI

In many of our experiments, utilising human-readable IRI-based representations linking to DBpedia provide best results. Our choice is motivated by IRIs representing the most straightforward – even if knowledge base-dependent – representation for entities. We assume for this to likely be the case due to these being succinct, human-readable, similar to labels and that the start of the candidate (<https://dbpedia.org/resource/>) provides a useful bias to contextually reinforce the representation. We provide an example IRI-based candidate representation prompt in Figure 5.

A.1.2 Label

A second candidate representation type we investigate is labels. The choice is motivated by prior work and structural similarity to human-readable Wikipedia- and DBpedia-based IRIs.

A.1.3 Description

While some of our other representations are likely to have occurred in appropriate contexts within chosen models’ training data, representing candidate entities as descriptions follows the rationale

SYSTEM

You are an expert assistant disambiguating entities and outputting if any of the passed entities are referenced in a given input text.

USER

Identify which entity candidate (if any) corresponds to the mention "Mika" in the input document.
Please reply with just the IRI of the entity.

Input document:

"Mika left his mark on Grand Prix history."

Entity Candidates:

- https://dbpedia.org/resource/Mika_Häkkinen
- [https://dbpedia.org/resource/Mika_\(singer\)](https://dbpedia.org/resource/Mika_(singer))
- https://dbpedia.org/resource/FC_Mika
- https://dbpedia.org/resource/R._Mika
- https://dbpedia.org/resource/FC_Lahti
- https://dbpedia.org/resource/Mika_Nakashima
- https://dbpedia.org/resource/Mika_Singh
- https://dbpedia.org/resource/Mika_Kaurismäki
- https://dbpedia.org/resource/Mika_Waltari
- https://dbpedia.org/resource/Mika_Väyrynen

ASSISTANT

The correct

disambiguated entity is

https://dbpedia.org/resource/Mika_Häkkinen

Figure 5: Prompt - IRI: Entity Candidates represented by their DBpedia-grounded entity IRIs.

of "contextual reasoning". Descriptions describe entities, give deeper context and use more specific language to define what an entity represents at its core. As such, utilising descriptions for our experiments tests whether (1) longer bits of text may bias models into different directions and (2) models can reason over contexts, reaching desired answers. Unfortunately, our results with descriptions tend to be worse than with IRIs or labels. We theorise that ACDC’s strictly constrained decoding causes issue with descriptions due to looking to reproduce the exact description while descriptive texts may oftentimes begin with generic ambiguous formulas. A possible improvement could be to change the decoding algorithm to beam search rather than the oftentimes defaulted-to greedy decoding. Further, descriptions may greatly vary in length, causing an underlying DLM to run into unexpected "behaviour" regarding attention and positional embeddings when a large "gap" is perceived between

Entity Representation	AGNUS w. IRI	AGNUS w. Type	AGNUS w. Label	AGNUS w. Desc.
AGNUS w. IRI	0.867	0.855	0.763	0.854
AGNUS w. Type	<u>0.855</u>	0.705	0.734	<u>0.766</u>
AGNUS w. Label	0.763	0.734	0.743	0.744
AGNUS w. Desc.	0.854	<u>0.766</u>	<u>0.744</u>	0.763

Table 4: Ablation Study (Entity Representation - Single and Pairwise): Disambiguation results (F1-measure) on AIDA for pairwise and singular (diagonal) entity representation information types for candidates on AGNUS (Mistral): entity IRI, entity type(s), entity label and entity description. Per column top-ranked score in **bold**, second-ranked underlined.

# Documents	Mentions	# Type-Consistent Docs.	Type Consist. (Mean)
888	15,314	331	46.60%

Table 5: Some data statistics for AIDA-Syn. Type-Consistency compares pre-transformation types of entities to post-transformation types of entities and checks overlap.

System	AIDA-Syn	AIDA	ASM-10	ASM-50	ASM-100
Babelify (Moro et al., 2014)	0.7503	0.6729	0.7660	0.7111	0.6912
WAT (Piccinno and Ferragina, 2014)	0.8641	0.6986	0.9355	0.8235	0.8332
REL (van Hulst et al., 2020)	-	-	0.9030	0.7942	0.6829

Table 6: F1 measures on datasets AIDA-Syn, AIDA for AGNUS and GERBIL-available systems (all other publicly available systems on GERBIL (Verborgh et al., 2018) timed out or returned "The annotator caused too many single errors" for the platform despite repeated attempts).

candidate descriptions' number of ingested tokens.

A.1.4 Type

Finally, we introduce 'types' as a source of information for entities. While simple and potentially ambiguous, the idea was that in combination with other representation types, it could help improve disambiguation by providing more context as used in more traditional entity disambiguation approaches (e.g. applying named entity recognition incl. types and disambiguating based on type-filtered candidates). On another hand, types can also be particularly specific, such as defining an entity as a "Formula One racer" which would prove beneficial to identify a mention Mika as the Finnish-born race driver Mika Häkkinen.

A.1.5 Pairwise Representation

In Table 4, we apply disambiguation based on multiple entity representations in a pairwise fashion AIDA. These experiments' prompts are set up analogously to the ones as illustrated in Figures 5 and 6 In these experiments, pairwise representations are ordered in descending fashion by mean representation scores reached in single-representation experiments (see Table 2): IRI > Label > Description > Type. For instance, in a pair of Label and Type, "candidate representa-

tion type 1" would be defined by Label and "candidate representation type 2" by Type.

A.2 Candidate Representation - Pairwise Effects

We investigated effects of single candidate representation types within our paper. We considered it interesting to have a look at pairwise combinations thereof as well to verify to what extent adding more information could yield better results – as would be an initial human intuition.

In Table 4, we evaluated AGNUS on pairwise combinations of candidate representation types to verify effects as well as the extent of increased information content on results. We note that disambiguating based on meaningful IRIs, such as from Wikipedia (e.g. [https://en.wikipedia.org/wiki/Mika_\(singer\)](https://en.wikipedia.org/wiki/Mika_(singer))), yields the best scores regardless of representation it may be combined with. Any further representation type worsens results, seemingly indicating that highly-defining compact representations may yield best results.

Types by themselves return mixed results, slightly improving upon description-based candidates, but deteriorate label-based results slightly. This may be due to the high overlap among candidates for this representation, potentially causing confusion upon disambiguation and yielding worst results (7.0%) in our experiments. Adding labels (7.3%) or descriptions (7.7%) to types increases candidate information, decreasing ambiguity and leading to improved results. Labels as an entity characteristic by themselves (7.43%) are relatively ambiguous, but benefit slightly from further information in the form of descriptions (7.44%). Overall, top scores are reached with IRI representations regardless of other combined information – actu-

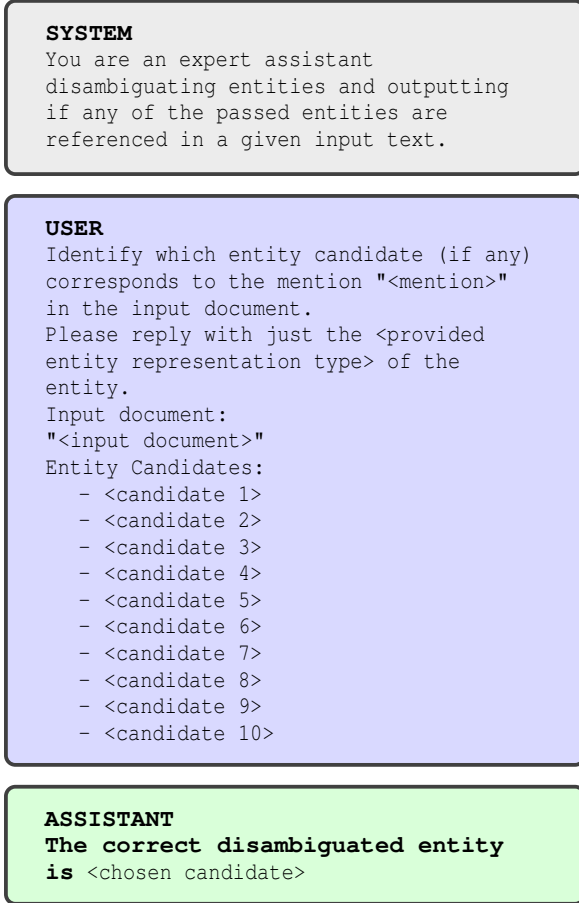


Figure 6: Prompt - Template: Template variables are surrounded by less than (<) and greater than (>) symbols.

ally suffering from any additional representations (by itself: 8.67%, with type(s): 8.55%, with description: 8.54%) –, most notably suffering from labels (7.63%).

A.3 Masked Attention Candidate Set - Details

MACS hides certain tokens' positions from other tokens without requiring retraining by restricting attention and sharing positional embeddings for successive token predictions. In Figure 8, we illustrate an example of the encoding to an underlying DLM: the positional embedding ID for each candidate is reset to the first candidate's positional embedding and incremented for each token until another candidate entity or the end of the candidate entity set is encountered. The positional embedding of the first token succeeding a masked attention candidate set is set to the longest candidate's final token's positional ID incremented by one.

As such, to the underlying DLM it will appear as though there was a gap in positions. Consequently,

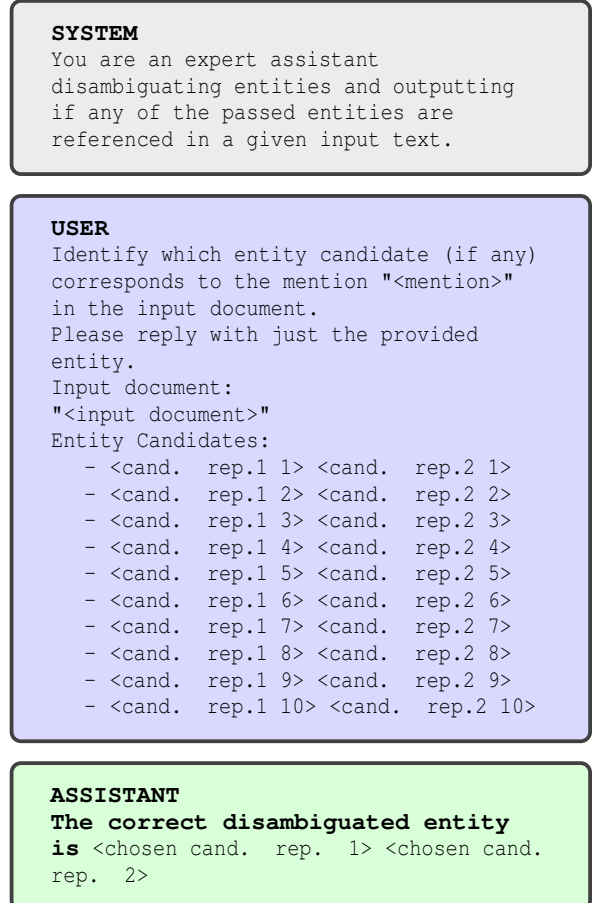


Figure 7: Prompt - Pairwise Template: Prompt Template for pairwise representation type experiments. Template variables are surrounded by less than (<) and greater than (>) symbols.

prompts including candidates with highly varying lengths may lead to weirdness for the underlying DLM's decoding process.

A.4 Agnus Contextual Decoding - Details

Our constrained decoding mechanism functions in a tree-based fashion in accordance to the grammar defined in Section 3.3 and only allows for specific tokens at each step by setting disallowed tokens to negative infinity ($-\infty$) with a customised logit processor. Effectively, this leads to undesired tokens being impossible to be generated by the language model. Our approach aims to be minimally invasive and maximally generalisable to open-weight models in the sense that we do not define nor modify a particular decoding algorithm. Instead, our models use the default or otherwise defined decoding algorithms for the respective language model (i.e. *greedy decoding*, *top-p sampling* or similar).

Furthermore, we introduce a simple optimization

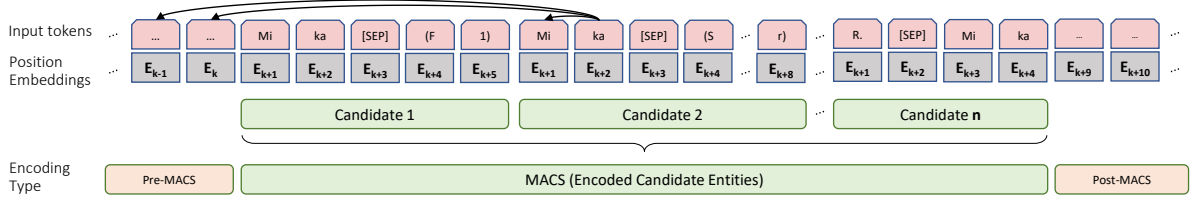


Figure 8: MACS Example - Position Embeddings & Attention: Encoding of a disambiguation prompt (see Fig. 6). In particular with attention from second candidate entity (Mika (Singer)) token <ka> with intra- and extra-candidate attention. Candidate entities attend (arrows) to themselves, but cannot attend to each other as defined in our causal mask (see Figure 3).

involving an "early stopping" mechanism when decoded entity representations start being no longer ambiguous

A.5 AIDA-Syn

We created AIDA-Syn using DeepSeek-R1 and generated 5 variants of coherent mentions and entities each. We automatically filtered out variants and documents where entities did not correspond to a valid DBpedia entity or where other LLM-related issues may have arisen. Some issues were related to DeepSeek’s maximum number of generated tokens activating prior to reaching the end. Key criteria for our generation included semantic coherence, lexical diversity, naturalness, plausibility within the surrounding text, and alignment with existing entities – our employed setups, prompts and results are publicly available¹². In the end, we generated a collection of 888 synthetic documents with 5 variants of mentions and entities for each. The reasoning behind the latter being more options for future evaluation endeavours, as well as switching to alternative mention contexts in case of faulty generations. Note that due to hardware limitations, we relied on DeepSeek’s API rather than employing ACDC with DBpedia entities for the synthetic dataset generation.

As a means of verifying that our generated mentions and entities are sensical, we used a two-pronged approach. First, two researchers manually validated a random sample of 222 (25%) documents, verifying contextual coherence for all variants. In 14% (32 documents) of documents, generated mentions created excessive entity ambiguity or were incorrect, leading to using another set of generated mentions and entities for the documents.

Second, we attempted to run the full suite of annotators via GERBIL (Verborgh et al., 2018) to see whether existing approaches could annotate

documents effectively – we report the results in Table 6. Unfortunately, many D2KB annotators did not run on our full AIDA-Syn (or ASM-10¹³, ASM-50¹⁴, ASM-100¹⁵) and the original AIDA datasets, instead returning timeout errors and similar. In Table 5 we display some details about the synthetic AIDA-Syn dataset including number of documents (888), total number of mentions (15,314) as well as entity type consistency between the original dataset and the transformed documents. Our assumption is that a certain degree of overlap between types should persist, but that it shouldn’t be an absolute overlap the sake of document diversity.

A.5.1 Synthetic Data Generation Caveats

Generating novel mentions can result in a variety of ways that lessen expressiveness of data. For instance, "Mika" is transformed into following alternative mentions: Ayrton Senna, Michael Schumacher, Alain Prost, Lewis Hamilton and Sebastian Vettel. In this case, a simple first name ("Mika") is transformed into a first name followed by a last name, both belonging to famous race drivers. The presence of the last name creates a lower degree of ambiguity than in the initial dataset, reducing complexity of the disambiguation task.

Further, in this setup, the linked wikipedia ID ("wiki") may be hallucinated, making the alternative document effectively unuseable unless verified.

Similarly, a DLM may hallucinate a mention and link an unrelated entity to it.

These points put into question the validity of utilising synthetic data for evaluation, but still allow for a certain degree of expressiveness regarding contamination diagnostics.

¹³<http://gerbil.aksw.org/gerbil/experiment?id=202505190000>

¹⁴<http://gerbil.aksw.org/gerbil/experiment?id=202505190001>

¹⁵<http://gerbil.aksw.org/gerbil/experiment?id=202505190002>

¹²<https://github.com/kmdn/agn-dis>

A.6 Baseline Experiments


Do note that in the case of "*w.o. ACDC*" (without constrained decoding), we apply fuzzy matching between candidate representations for both predicted and expected values, ranking similarity for the sake of comparison fairness and picking the highest-overlap-similarity candidate as a match. Just using the results as-is for a "baseline" comparison seemed disingenuous as "exact matching" criteria would put baseline results very close (if not exactly) to 0 in most cases.

Applying hard-prompting based finetuning to our employed suite of large language models would likely alleviate the effects to a certain degree, but would simultaneously render the comparison invalid due to comparing our zero-shot model to a 1-shot baseline, therewith having only limited expressivity over our existing ED evaluation table (Table 1).

Due to similar reasons, our baseline without candidates still uses matching to candidates (it did not see or produce) rather than dryly applying an exact matching scheme, therewith heightening the likelihood of correct results. Hence, we urge readers to not overestimate baseline performance.

A.7 Evaluating with Related Work Candidates & Details

Unfortunately, to the best of our knowledge, the large majority of prior work does not provide candidate entities for their entity disambiguation methods. We provide our candidates, data and results (see code repository). We have some comments and concerns with prior work's (EntGPT (Ding et al., 2024a)) provided candidates¹⁶, but we want to validate our approach best as possible and see an added benefit for comparability in doing so.

AGNUS  always chooses a single entity from among a set of 10 candidates. EntGPT's candidate sets are variable in size and less than or equal in amount to 10 (averaging between 6.1 - 9.2 depending on dataset, see Table 7). Further, EntGPT relies on a prompt structure allowing underlying DLMs to express that "None" of the provided candidates correspond to the desired one. This unfortunately entails a few additional considerations regarding evaluation: Theoretically, (1) if a candidate generation technique were to never contain a desired

¹⁶https://github.com/yifding/In_Context_EL/tree/main/RUN_FILES/4_13_2023/rel_blink/evaluation_new_one_step

Dataset	In-Set	"None"	Docs.	Docs. (Original)	Mentions	Total Cand.	Min. Cand.	Avg. Cand.
KORE50	113	35	50	50	148	1365	1	9.223
ACE2004	242	15	35	106	257	1953	1	7.599
AIDA-B	4250	125	230	231	4375	31651	1	7.235
AQUAINT	700	27	50	50	727	4935	1	6.788
CLUEWEB	9961	1193	320	-	11154	83285	1	7.467
MSNBC	617	39	20	20	656	4525	1	6.898
OKE2015	441	95	101	101	536	3625	1	6.763
OKE2016	240	48	55	-	288	2179	1	7.566
Reuters-128	544	106	113	128	650	4686	1	7.209
RSS-500	447	77	357	500	524	3199	1	6.105
WIKI	6076	717	319	-	6793	42296	1	6.226

Table 7: Prior Work (Ding et al., 2024a) Dataset Statistics: Number of mentions for which correct entity is within candidate set (In-Set), is not in candidate set ("None" being correct), number of documents provided and number of documents within dataset originally (as far as could be determined reasonably). Dash (-) means varying values have been found from different sources.

Dataset	In-Set	"None"	Mentions	Total Cand.	Min. Cand.	Max. Cand.	Avg. Cand.
ACE2004	2333	274	2607	257718	1	100	98.86
AQUAINT	13359	1383	14742	1464034	1	100	99.31
AIDA	25076	2741	27817	2760855	0	100	99.25
KORE 50	813	93	906	89440	6	100	98.72
MSNBC	10450	1066	11516	1143092	1	100	99.26

Table 8: AGNUS Dataset Statistics: Candidates generated using DBpedia Lookup. Number of mentions for which correct entity is within candidate set (In-Set), is not in candidate set ("None" being correct), number of mentions, number of total candidates along with minimum, maximum and average candidates for all mentions.

entity for disambiguation, a DLM could technically always choose "None" and reach a perfect score. (2) Limited comparability to existing ED methods.

Also, according to our analyses (see Table 7), some datasets are incomplete in terms of documents and mentions, therefore making meaningful comparisons with other existing work difficult. Despite concerns regarding generalizability to other methods, we regard comparing AGNUS to the best of our knowledge the only prior work that explicitly provides candidate entities a meaningful endeavour. We note that while some prior work do technically provide code to generate candidates, provided code being impossible to run without possible major changes (e.g. local dependencies (Liu et al., 2024)¹⁷ or paths to inaccessible datasets (Xiao et al., 2023a)¹⁸) significantly impedes a comparable and clean evaluation process.

¹⁷https://github.com/laquabe/OneNet/blob/main/pointwise_process/listwise_cand.py

¹⁸https://github.com/MrZilinXiao/InsGenEntityLinking/blob/master/data_scripts/create_candidates_dict.py